

Algorithm 1 OPNP - Optimal Parameter and Neuron Pruning

Require: A trained model $f(\mathbf{x}; \boldsymbol{\theta})$, training samples $\{\mathbf{x}_i\}_{i=1}^m$, pruning percentage $\rho_{max}^w, \rho_{min}^w, \rho_{min}^o, \rho_{max}^o$, test sample $\{\mathbf{x}_k\}_{k=1}^n$;

Ensure: Energy score of test sample $\{E(\mathbf{x}_k)\}_{k=1}^n$

- 1: **for** $i = 1$ to m **do**
- 2: compute gradient $g(\mathbf{x}_i)$;
- 3: **end for**
- 4: $\mathbf{M} = \frac{1}{m} \sum_{k=1}^m |g(\mathbf{x}_k)|$;
- 5: $\mathbf{O}_i = \frac{1}{K} \sum_{p=1}^K \mathbf{M}_{ip}, \quad i = 1, \dots, L$;
- 6: Get threshold $\Omega_{min}^w, \Omega_{max}^w$ by ranking \mathbf{M} ;
- 7: Get threshold $\Omega_{min}^o, \Omega_{max}^o$ by ranking \mathbf{O} ;
- 8: $\mathbf{W}[\mathbf{M} > \Omega_{max}^w] = 0$ and $\mathbf{W}[\mathbf{M} < \Omega_{min}^w] = 0$;
- 9: **for** $k = 1$ to n **do**
- 10: Compute the pre-logit embedding $h(\mathbf{x}_k)$;
- 11: $h(\mathbf{x}_k)[\mathbf{O} > \Omega_{max}^o] = 0$ and $h(\mathbf{x}_k)[\mathbf{O} < \Omega_{min}^o] = 0$;
- 12: $f(\mathbf{x}_k) = \mathbf{W} \cdot h(\mathbf{x}_k) + \mathbf{b}$;
- 13: $E(\mathbf{x}_k) = -\log \sum_{i=1}^K \exp(f_i(\mathbf{x}_k))$;
- 14: **end for**
- 15: **Return** $\{E(\mathbf{x}_k)\}_{k=1}^n$

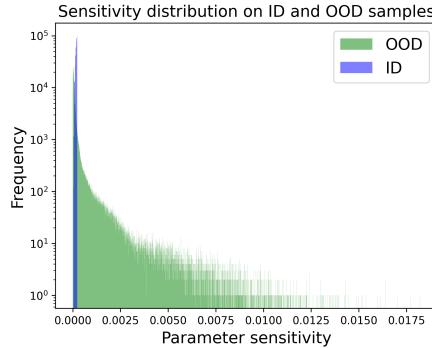


Figure 1: Sensitivity distribution on ID and OOD samples. The top 20% sensitive parameters (measure on ID set) are illustrated. The average sensitivity on OOD samples is 0.00024 and the average sensitivity on ID samples is 0.00018.

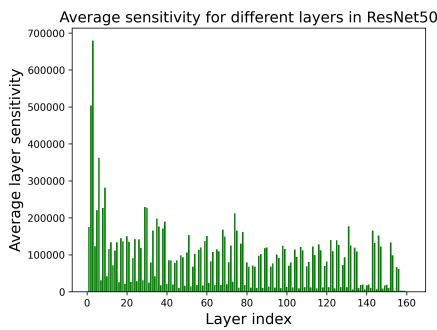


Figure 2: Average parameter sensitivity for different layers in ResNet50.